# Auto-tuning shared memory parallel Sparse BLAS operations in `librsb`-1.2

**Michele MARTONE**   (michele.martone@ipp.mpg.de)

Max Planck Institute for Plasma Physics, Garching bei München, Germany

## Intro

- Sparse BLAS (Basic Linear Algebra Subroutines) [2] specifies main kernels for iterative methods:
  - *sparse Multiply by Matrix*: "**MM**"
  - *sparse triangular Solve by Matrix*: "**SM**"
- Focus on **MM**: $C \leftarrow C + \alpha\, op(A) \times B$, with
  - $A$ has dimensions $m \times k$ and is *sparse* ($nnz$ nonzeroes)
  - $op(A)$ can be either of $\{A, A^T, A^H\}$ (parameter $transA$)
  - *left hand side* (**LHS**) $B$ is $k \times n$,
    *right hand side* (**RHS**) $C$ is $n \times m$ ($n$=NRHS),
    both dense (eventually with *strided* access $incB, incC, ...$)
  - $\alpha$ is scalar
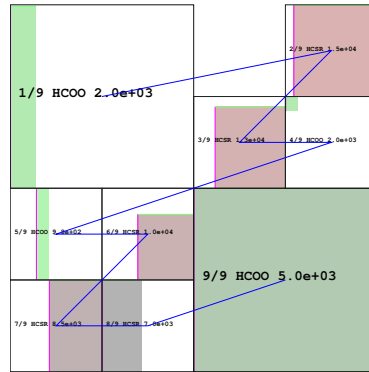  - either single or double precision, either `real` or `complex`
- `librsb` implements the Sparse BLAS using the **RSB** (**Recursive Sparse Blocks**) data structure [3].
- Hand tuning for each operation variant is impossible.
- We propose **empirical auto-tuning** for `librsb`-1.2.
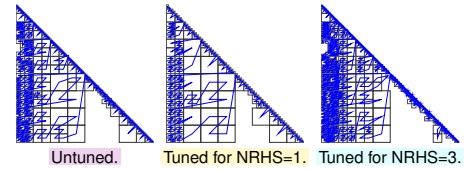
## RSB: Recursive Sparse Blocks

- *Sparse blocks* in COO or CSR [1].
- ...eventually with 16-bit indices ("HCOO" or "HCSR").
- *cache blocks* suitable for *thread parallelism*.
- *Recursive partitioning* of submatrices results in Z-ordered blocks.



- Instance of matrix *bayer02* (ca. $14k \times 14k$, $64k$ nonzeroes).
- The black-bordered boxes are *sparse blocks*.
- Greener have fewer nnz than average, redder have more.
- Blocks rows (columns) of LHS (RHS) range during **MM**.
- Larger submatrices like "**9/9**" can have fewer nonzeroes than smaller ones like "**4/9**".

## Merge / split based autotuning

- Optimal default blocking ?
  - Irregular matrix patterns !
  - Operands (especially $transA$, NRHS) change memory footprint !
- Empirical auto-tuning:
  - Given a Sparse BLAS operation, probe for a *better performing* blocking.
  - Search among *slightly coarser* or *finer* ones.



Untuned.   Tuned for NRHS=1.   Tuned for NRHS=3.

- Tuning example on symmetric matrix *audikw_1*.
- Here only lower triangle, ca. $1M \times 1M$, $39M$ nonzeroes.
- On a machine with 256 KB sized L2 cache.
- Left one (625 blocks, avg 491 KB) is before tuning.
- Middle one (271 blocks, avg 1133 KB) after tuning (1.057x speedup, 6436.6 ops to amortize) for **MV** (**MM** with NRHS=1).
- Right one (1319 blocks, avg 233 KB) after tuning (1.050x speedup, 3996.2 ops to amortize) for **MM** with NRHS=3.
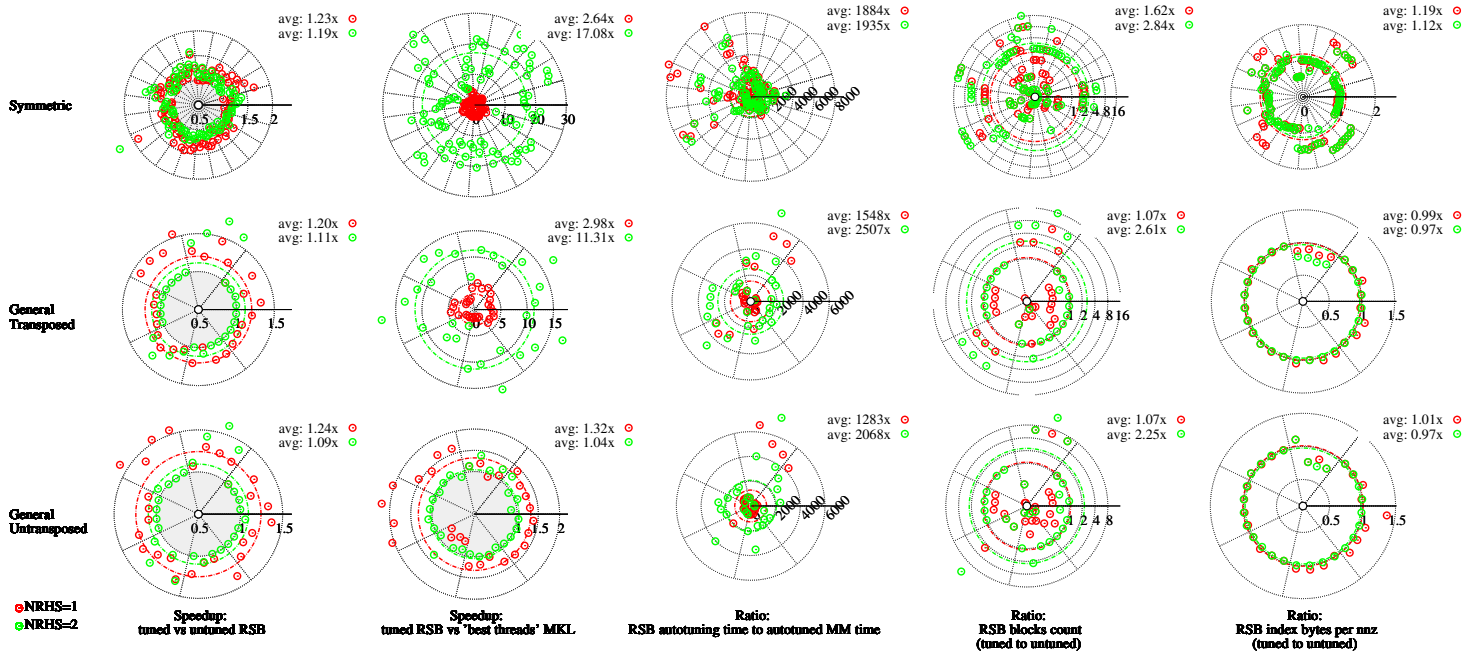- Finer subdivision at NRHS=3 consequence of increased cache occupation of per-block LHS/RHS operands.

### Sparse BLAS autotuning extension

```
1  ! Matrix-Vector Multiply: y ← alpha*op(A)*x+y
2  call USMV(transA,alpha,A,x,incx,y,incy,istat)
3  ! Tuning request for the next operation
4  call USSP(A,blas_autotune_next_operation,istat)
5  ! Matrix structure and threads tuning
6  call USMV(transA,alpha,A,x,incx,y,incy,istat)
7  ...
8  do ! A is now tuned for y ← alpha*op(A)*x+y
9     call USMV(transA,alpha,A,x,incx,y,incy,istat)
10 ...
11 ! Request autotuning again
12 call USSP(A,blas_autotune_next_operation,istat)
13 ! Now tune for C ← C + alpha * op(A) * B
14 call USMM(order,transA,nrhs,alpha,A,B,ldB,C,ldC,istat)
15 ! The RSB representation of A is probably different than before USMM
```
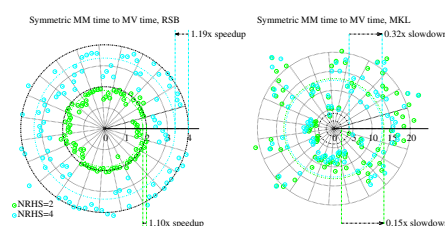
## Experiment in MM tuning and comparison to MKL



NRHS=1
NRHS=2

Symmetric — avg: 1.23x, avg: 1.19x / avg: 2.64x, avg: 17.08x / avg: 1884x, avg: 1935x / avg: 1.62x, avg: 2.84x / avg: 1.19x, avg: 1.12x

General Transposed — avg: 1.20x, avg: 1.11x / avg: 2.98x, avg: 11.31x / avg: 1548x, avg: 2507x / avg: 1.07x, avg: 2.61x / avg: 0.99x, avg: 0.97x

General Untransposed — avg: 1.24x, avg: 1.09x / avg: 1.32x, avg: 1.04x / avg: 1283x, avg: 2068x / avg: 1.07x, avg: 2.25x / avg: 1.01x, avg: 0.97x

Speedup: tuned vs untuned RSB | Speedup: tuned RSB vs 'best threads' MKL | Ratio: RSB autotuning time to autotuned MM time | Ratio: RSB blocks count (tuned to untuned) | Ratio: RSB index bytes per nnz (tuned to untuned)

## Setup

- `librsb` (`icc -O3 -xAVX`, v15) vs Intel MKL (v.11.2) CSR.
- $2\times$ "*Intel Xeon E5-2680*", 16 OpenMP threads.
- **MM** with NRHS={**1,2**}, four BLAS numerical types.
- 27 matrices in total (as in [3]), including symmetric.

## Results Summary

- **Few dozen percent** improvement over untuned, costing **few thousand** operations.
- **Significantly faster than Intel MKL** on symmetric and transposed operation with NRHS=2 ($> 1$).
- **Autotuning** more effective on **symmetric** and **unsymmetric untransposed** with **NRHS=1**.
- Tuning mostly **subdivided further** for NRHS=2.

## Highlight: symmetric MM vs MV performance



Symmetric MM time to MV time, RSB — 1.19x speedup / 1.10x speedup
Symmetric MM time to MV time, MKL — 0.32x slowdown / 0.15x slowdown

NRHS=2
NRHS=4

RSB Symmetric **MM** performance increases when NRHS>1, while Intel MKL CSR' falls. See here for NRHS=2 and NRHS=4.

## Outlook

One may improve via:

- **Reversible in-place merge and split**: no need for copy while tuning.
- Best merge/split choice not obvious: **different merge and split rules**.
- **Non-time efficiency criteria** (e.g. use an energy measuring API when picking **better performing**).

## References

[1] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. V. der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods, 2nd Edition*. SIAM, Philadelphia, PA, 1994.

[2] BLAS Forum. Sparse BLAS specification (BLAS specification, chapter 3). Technical report.

[3] M. Martone. Efficient multithreaded untransposed, transposed or symmetric sparse matrix-vector multiplication with the recursive sparse blocks format. *Parallel Computing*, (40):251–270, 2014.